

Speech Emotion Recognition and Perception of Music

Mélanie Fernández Pradier

Prof. Dr.-Ing. Bin Yang

Supervisors: Prof. Dr.-Ing. Bin Yang
Dipl.-Ing. Fabian Schmieder

January 27, 2011

Motivation

Speech Emotion Recognition and Perception of Music

Emotion Recognition from Speech

- Speech \sim two-channel
 - linguistic
 - paralinguistic
- Several Applications
 - support ASR
 - diagnoses
 - speech synthesis
 - entertainment

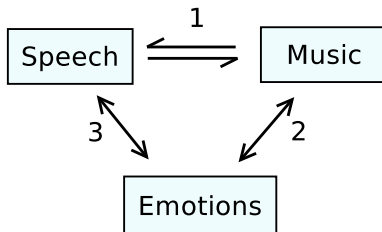
Music Perception

- “language of emotion”
- treatment of affective disorders
- treatment of speech disorders
- same origin of music and speech

Aim of the thesis

Apply Music Theory to Speech Emotion Recognition

Investigate Speech and Music similarities to derive universal features for Emotions



- ① What is the link between music and speech?
- ② How are emotions transmitted through music?
- ③ Can we apply musical knowledge to speech processing?

1 Introduction

- Motivation
- Aim of the thesis

2 Basic Features

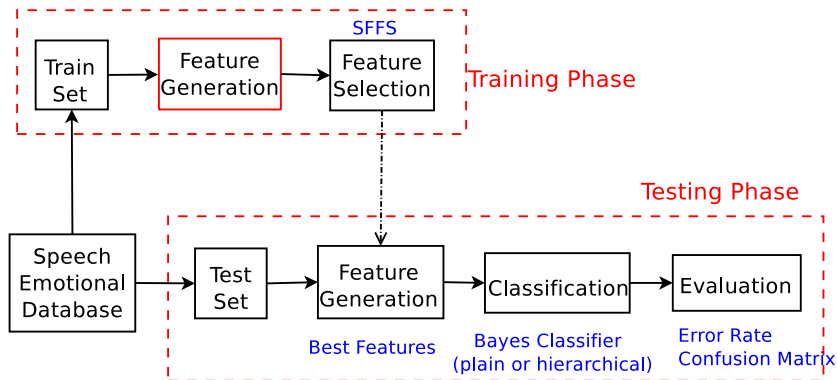
- General Concepts
- Description

3 Musical Features

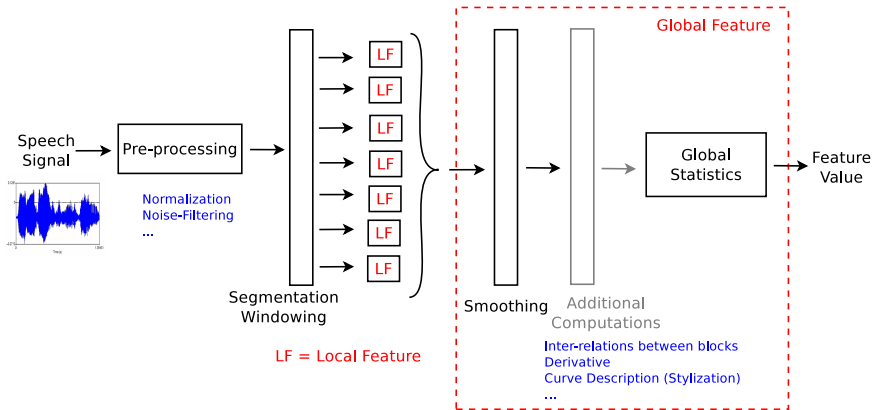
- Interval and Triad Features
- Based on Music Emotion Recognition
- Perceptual Model of Intonation

4 Simulations and Results

Pattern Recognition



Feature Generation



Basic Features Description

Local Features

- ZCR
- MFCC
- Energy
total + bands
- Pitch
- Voiced-unvoiced
- VAD

$$ZCR = \frac{1}{2} \cdot \sum_{n=1}^N |\text{sgn}(x_n) - \text{sgn}(x_{n+1})|$$

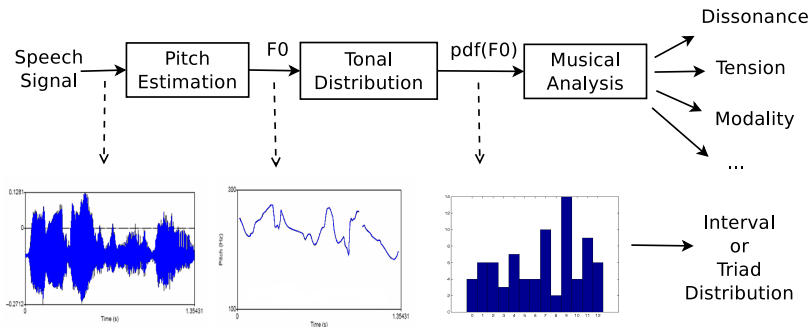
$$Cepstrum = \left| FFT \left\{ \log \left(|FFT \{ \underline{x} \}|^2 \right) \right\} \right|^2$$

$$Energy = \sum_{n=1}^N x_n \cdot x_n^*$$

Global Features

- Global statistics: min, mean, max, median, std, iqr...
directly, 1st or 2nd derivative
- Energy and pitch plateaux
- Combination with logical features

Interval and Triad Features



Interval Features

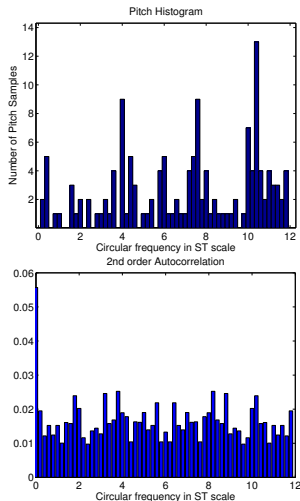
Autocorrelation of the circular pitch density function

$$\int_0^L p_o(\text{mod}_L(s + \lambda)) p_o(\lambda) d\lambda$$

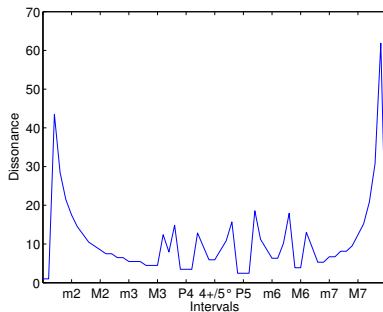
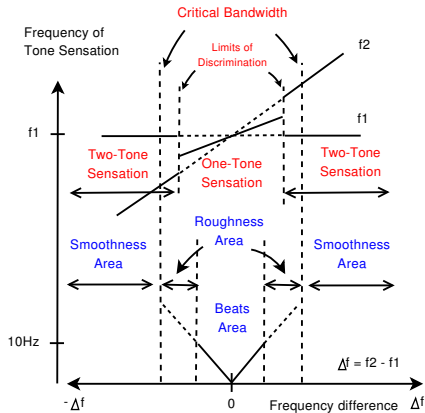
Intervalic dissonance

$$DIS = \int_0^L d(s) r_o(s) ds$$

$$\text{where } d(s) \simeq \sqrt{N(s) D(s)}$$

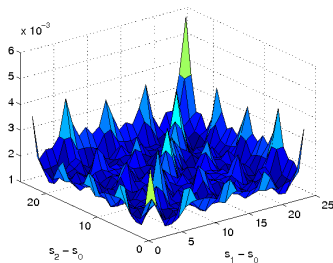


Interval Dissonance

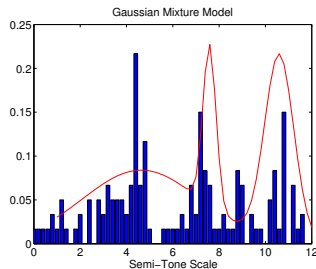


Triad Features

- 1 Direct computation
- 2 Extraction of “dominant pitches”



Autocorrelation Triad Features



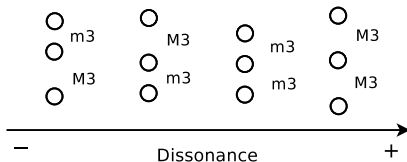
Gaussian Triad Features

Tension and Modality



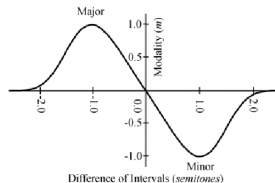
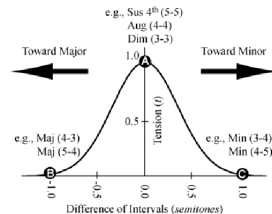
Stable, final
resolved

Unstable, tense
unresolved



m3 = minor third (3 ST)

M3 = major third (4 ST)



Loudness, Timbre and Rhythm

Intensity Features

$$I(k) = \sum_{n=0}^{N/2} |\text{FFT}_k(n)|$$

$$D_i(k) = \frac{1}{I(k)} \sum_{n=L_i}^{H_i} |\text{FFT}_k(n)|$$

where k refers to the frame

Timbre Features

$$\text{FFT}_k \equiv \{x_{k1} \dots x_{kN}\}$$

$$\rightsquigarrow \text{sorted} \equiv \{x'_{k1} \dots x'_{kN}\}$$

$$\text{Peak}(k) = \log \left\{ \frac{1}{\alpha N} \sum_{i=1}^{\alpha N} x'_{ki} \right\}$$

$$\text{Valley}(k) = \log \left\{ \frac{1}{\alpha N} \sum_{i=1}^{\alpha N} x'_{k(N-i+1)} \right\}$$

Loudness, Timbre and Rhythm

Rhythm Features

- 1 Compute FFT
- 2 Extract amplitude envelope

$$A_i(n) = FFT_i(n) \otimes h_w(n)$$

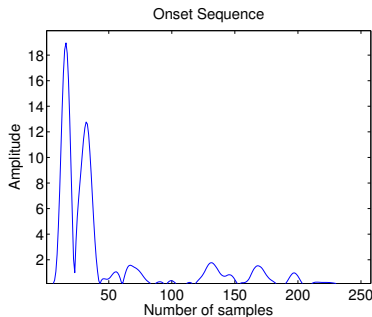
- 3 Apply Canny operator

$$O_i(n) = A_i(n) \otimes C(n)$$

$$C(n) = \frac{n}{\sigma^2} e^{-\frac{n^2}{2\sigma^2}}$$

We obtain the onset sequence

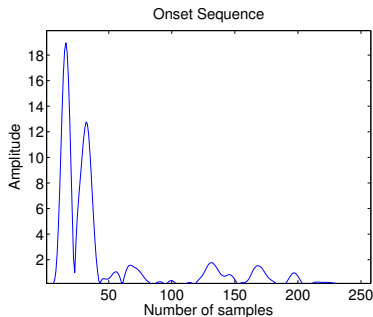
$$O_i(n)$$



Loudness, Timbre and Rhythm

Rhythm Features

- Strength** Average value of the peaks
- Regularity** Average value of peaks in the autocorrelation
- Speed** Ratio of number of peaks and time duration



Perceptual Model of Intonation

Perceptual principles

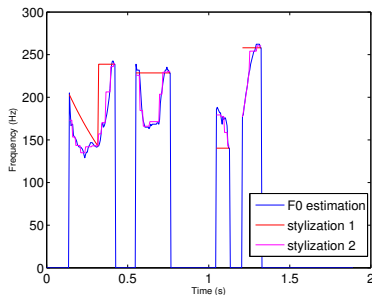
- 1 Segmentation Effect
- 2 Glissando Threshold: minimum amount of frequency change

$$g_{th} = 0.16 / T^2 \text{ [ST/s}^2\text{]}$$

- 3 Differential Glissando Threshold: minimum difference in slope

$$dg_{th} = a_2 - a_1 = 20 \text{ [ST/s]}$$

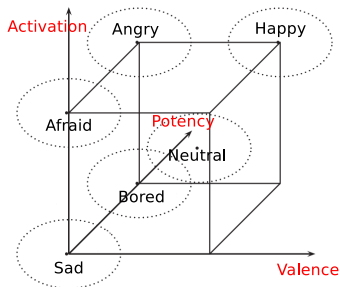
- 4 Short-term integration in time



Database - Labels - Features

Database: emoDB (TUB)

- 10 speakers
- 708 files
- 6 emotions

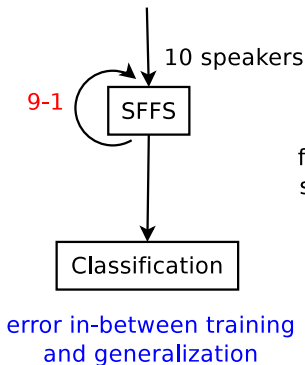


BASIC	SET
duration	16
MFCC	91
ZCR	13
harmony	3
energy	58
pitch	33
Total	214

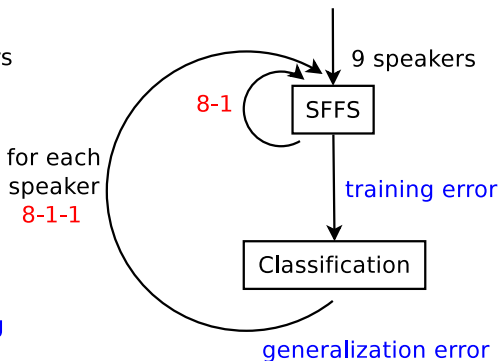
MUSICAL	SET
interval	31
autocorr. triad	4
gaussian triad	10
intensity	63
rhythm	15
Total	123

Strategies for evaluation 9-1 Vs 8-1-1

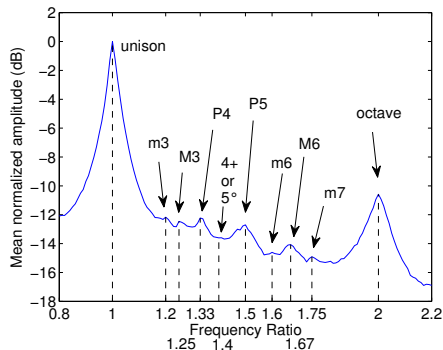
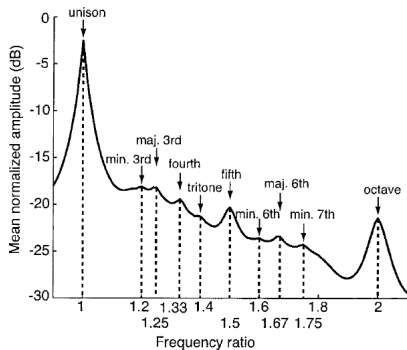
Evaluation 9-1



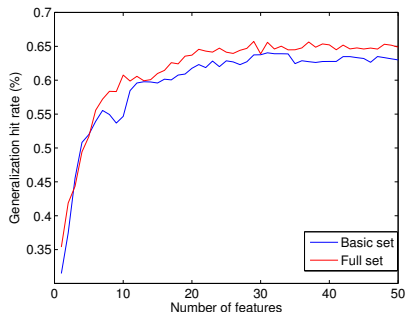
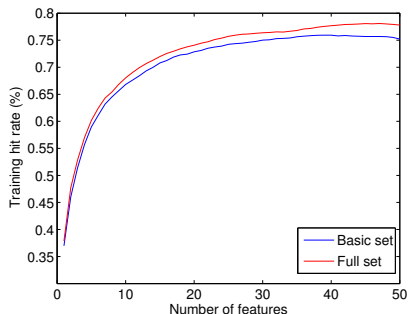
Evaluation 8-1-1



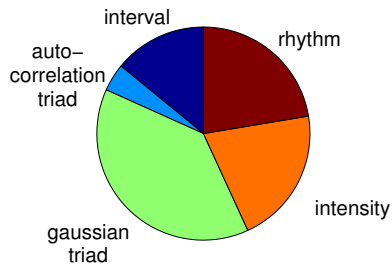
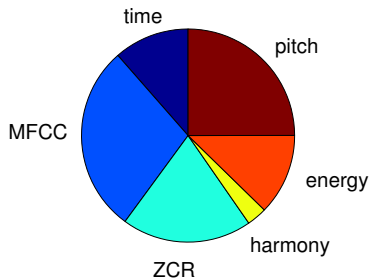
Musical Universals



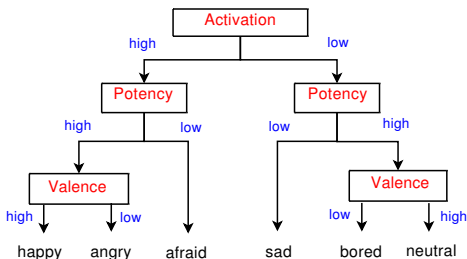
Plain bayes classifier - Evaluation 8-1-1



Nature of selected features

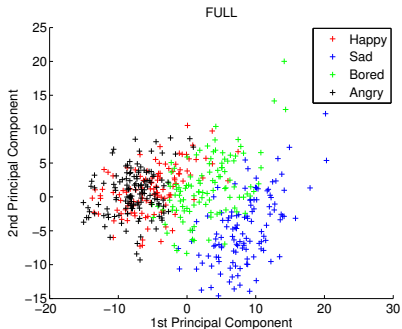
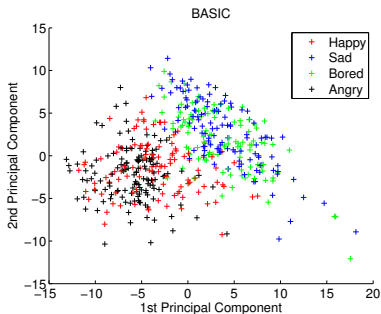


Comparison plain Vs hierarchical bayes classifier

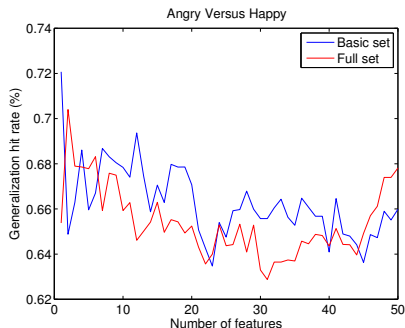
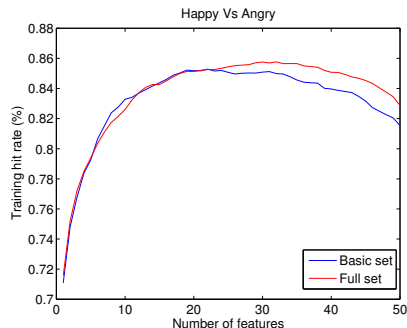


	plain Bayes	hierarchical Bayes
Basic	76.12	84.22
Basic+ Interval+Triad	80.61	85.04

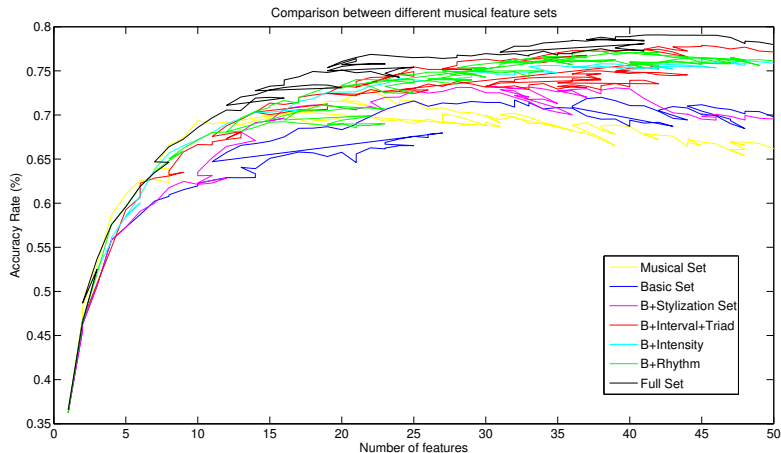
Multi-dimensional Scaling



Happy Vs Angry - Evaluation 8-1-1



Final Comparison of Musical Features



Conclusion

Summary

- ① Literature review about speech, music and emotions
- ② Theoretical background on psychoacoustics
- ③ Re-implementation of the basic features
- ④ Implementation of speech processing algorithms
- ⑤ Implementation of musical features
(music perception, MER and linguistics)
- ⑥ Simulations \Rightarrow Musical features can help to improve emotion recognition in speech

Conclusions

Further research

- Environment: natural emotional speech, other languages
- Pattern Recognition steps: feature transformation, pitch extraction, classification...
- Improvement of musical features
Dissonance model, Perceptual model of intonation,
Emotionally meaningful moments
- Systematization of feature extraction step

"Even monkeys express strong feelings in different tones – anger and impatience by low, – fear and pain by high notes."

Charles Darwin, Naturalist

Thank you!

Looking forward to your questions. . .